

WHAT IS CLAIMED IS:

- 1 1. A method of downloading data sets from among a plurality of host computers,
2 comprising the steps of:
3 (a) storing representations of data set addresses in a set of data structures,
4 including a buffer and a first disk file, wherein representations of data set addresses stored in
5 the first disk file are ordered;
6 (b) downloading at least one data set that includes addresses of one or more
7 referred data sets;
8 (c) identifying the addresses of the one or more referred data sets;
9 (d) for each identified address:
10 (d1) generating a representation of the identified address;
11 (d2) determining whether the representation is stored in the buffer, and when
12 this determination is negative, storing the representation in the buffer; and
13 (e) when the buffer reaches a predefined full condition:
14 (e1) ordering the contents of the buffer according to the representations; and
15 (e2) performing an ordered merge of the contents of the buffer into the
16 contents of the first disk file.
- 1 2. The method of claim 1, further comprising:
2 in step (d2), when the determination is negative, storing the identified address in the
3 buffer.
- 1 3. The method of claim 1, further comprising:
2 in step (d2), when the determination is negative, storing the identified address in a
3 second disk file;
4 in step (d2), additionally storing with each representation in the buffer a pointer to the
5 corresponding address stored in the second disk file; and
6 in step (e1), while ordering the contents of the buffer, keeping with each
7 representation in the buffer its pointer to the corresponding address in the second disk file.

1 4. The method of claim 3 wherein

2 step (e2) includes: for each representation in the buffer storing an associated flag,
3 setting the flag to a first value when the representation is equal to a representation previously
4 stored in the first disk file, and setting the flag to a second value, distinct from the first value,
5 when the representation is not equal to any representation previously stored in the first disk
6 file; and

7 step (e) includes: (e3) for each representation whose flag is set to the second value,
8 scheduling the corresponding data set for downloading.

1 5. The method of claim 1 wherein:

2 step (a), storing representations of data set addresses, includes the step of storing
3 representations of data set addresses in a sparse disk file which is divided into portions, each
4 portion having a starting address and contents comprising an ordered list of representations of
5 data addresses; and

6 step (e2), merging the contents of the buffer with the ordered contents of the sparse
7 disk file, includes:

8 for each of a plurality of the representations stored in the buffer:

9 (e2-1) determining a starting address for a corresponding portion of the sparse
10 disk file; and

11 (e2-2) performing an ordered merge of a subset of the buffer, starting at the
12 representation for which the starting address was obtained, into the contents of the
13 corresponding portion.

1 6. The method of claim 1 wherein:

2 step (a), storing representations of data set addresses, includes the step of storing
3 representations of data set addresses in a sparse disk file having empty entries interspersed
4 among entries storing said representations; and

5 step (e2), merging the contents of the buffer with the ordered contents of the sparse
6 disk file, includes:

7 for each respective representation stored in the buffer:

8 (e2-1) determining a starting address for a corresponding portion of the sparse

9 disk file; and

10 (e2-2) sequentially scanning the disk file, starting at the representation for
11 which the starting address was obtained, until the first of (A) a representation matching the
12 respective representation is found and (B) one of the empty entries is found, and when an
13 empty entry is found storing the respective representation in the empty entry.

1 7. The method of claim 1 wherein, in step (d1), the representation comprises a checksum
2 of at least a portion of the identified address.

1 8. The method of claim 1 wherein step (d2) further comprises:

2 (d2-1) determining whether the representation is stored in a cache before determining
3 whether the representation is stored in the buffer;

4 (d2-2) when the representation is not stored in the cache, the cache has not reached a
5 predefined full condition, and other predefined criteria are met, adding the representation to
6 the cache; and

7 (d2-3) when the representation is not stored in the cache, the cache has reached said
8 predefined full condition, and said other predefined criteria are met, evicting a stored
9 representation from the cache in accordance with an eviction policy and adding the
10 representation to the cache.

1 9. The method of claim 1 wherein step (e2) further comprises:

2 when a representation in the first buffer is not found in the first disk file during
3 merging, scheduling the corresponding data set for downloading.

1 10. The method of claim 8 wherein step (e2) further comprises:

2 when a representation in the buffer is not found in the first disk file during merging,
3 scheduling the corresponding data set for downloading.

1 11. The method of claim 8 wherein:

2 step (a), storing representations of data set addresses, includes the step of storing
3 representations of data set addresses in a sparse disk file which is divided into portions, each

4 portion having a starting address and contents comprising an ordered list of representations of
5 data addresses; and

6 step (e2), performing an ordered merge of the contents of the buffer into the contents
7 of the sparse disk file, includes:

8 for each of a plurality of the representations stored in the buffer:

9 (e2-1) obtaining a starting address for a corresponding portion of the sparse
10 disk file; and

11 (e2-2) performing an ordered merge of a subset of the buffer, starting at the
12 representation for which the starting address was obtained, into the contents of the
13 corresponding portion.

1 12. The method of claim 8 wherein:

2 step (a), storing representations of data set addresses, includes the step of storing
3 representations of data set addresses in a sparse disk file having empty entries interspersed
4 among entries storing said representations; and

5 step (e2), merging the contents of the buffer with the ordered contents of the sparse
6 disk file, includes:

7 for each respective representation stored in the buffer:

8 (e2-1) determining a starting address for a corresponding portion of the sparse
9 disk file; and

10 (e2-2) sequentially scanning the disk file, starting at the representation for
11 which the starting address was obtained, until the first of (A) a representation matching the
12 respective representation is found and (B) one of the empty entries is found, and when an
13 empty entry is found storing the respective representation in the empty entry.

1 13. A method of downloading data sets from among a plurality of host computers,
2 comprising the steps of

3 (a) storing representations of data set addresses in a set of data structures,
4 including a first buffer, a second buffer, and a first disk file, wherein the first disk file
5 contains ordered representations of data set addresses;

6 (b) selecting as a current buffer one of the first and second buffers;

7 (c) downloading at least one data set that includes addresses of one or more
8 referred data sets;
9 (d) identifying the addresses of the one or more referred data sets; and
10 (e) for each identified address:
11 (e1) generating a representation of the identified address; and
12 (e2) determining whether the representation is stored in the current buffer, and
13 when this determination is negative, storing the representation in the current buffer; and
14 (f) when the current buffer reaches a predefined full condition:
15 (f1) selecting the other buffer as the current buffer, wherein the previously
16 current buffer is identified as a non-current buffer;
17 (f2) ordering the representations stored in the non-current buffer; and
18 (f3) performing an ordered merge of the contents of the non-current buffer
19 into the contents of the first disk file.

1 14. The method of claim 13, further comprising:
2 in step (e2), when the determination is negative, storing the identified address in the
3 current buffer.

1 15. The method of claim 13, further comprising:
2 in step (e2), when the determination is negative, storing the identified address in a
3 second disk file;
4 in step (e2), additionally storing with each representation in the current buffer a
5 pointer to the corresponding address stored in the second disk file; and
6 in step (f2), while ordering the contents of the non-current buffer, keeping with each
7 representation in the non-current buffer its pointer to the corresponding address in the second
8 disk file.

1 16. The method of claim 15 wherein
2 step (e2) comprises: for each representation in the buffer storing an associated flag,
3 setting the flag to a first value when the representation is equal to a representation previously
4 stored in the first disk file, and setting the flag to a second value, distinct from the first value,

5 when the representation is not equal to any representation previously stored in the first disk
6 file; and

7 step (f) includes: (f4) for each representation whose flag is set to the second value,
8 scheduling the corresponding data set for downloading.

1 17. The method of claim 13 wherein step (e2) further comprises:

2 when a representation in the current buffer is not found in the first disk file during
3 merging, scheduling the corresponding data set for downloading.

1 18. The method of claim 13 wherein:

2 step (a), storing representations of data set addresses, includes storing representations
3 of data set addresses in a sparse disk file which is divided into portions, each portion having a
4 starting address and contents comprising an ordered list of representations of data addresses;
5 and

6 step (e2), performing an ordered merge of the contents of the current buffer into the
7 contents of the sparse disk file, comprises the following steps:

8 for each of a plurality of the representations stored in the current buffer:

9 (e2-1) obtaining a starting address for a corresponding portion of the sparse
10 disk file; and

11 (e2-2) performing an ordered merge of a subset of the current buffer, starting
12 at the representation for which the starting address was obtained, into the contents of the
13 corresponding portion.

1 19. The method of claim 13 wherein:

2 step (a), storing representations of data set addresses, includes the step of storing
3 representations of data set addresses in a sparse disk file having empty entries interspersed
4 among entries storing said representations; and

5 step (e2), merging the contents of the buffer with the ordered contents of the sparse
6 disk file, includes:

7 for each respective representation stored in the buffer:

8 (e2-1) determining a starting address for a corresponding portion of the sparse
9 disk file; and

10 (e2-2) sequentially scanning the disk file, starting at the representation for
11 which the starting address was obtained, until the first of (A) a representation matching the
12 respective representation is found and (B) one of the empty entries is found, and when an
13 empty entry is found storing the respective representation in the empty entry.

1 20. The method of claim 13 wherein, in step (e1), the representation comprises a
2 checksum of at least a portion of the identified address.

1 21. The method of claim 13 wherein step (e2) further comprises:

2 (e2-1) determining whether the representation is stored in a cache before determining
3 whether the representation is stored in the current buffer;

4 (e2-2) when the representation is not stored in the cache, and the cache has not
5 reached a predefined full condition, adding the representation to the cache; and

6 (e2-3) when the representation is not stored in the cache, and the cache has reached
7 said predefined full condition, evicting a stored representation from the cache in accordance
8 with an eviction policy and adding the representation to the cache.

1 22. A method of downloading data sets from among a plurality of host computers,
2 comprising the steps of:

3 (a) storing representations of data set addresses in a set of data structures,
4 including a buffer and a disk file, wherein representations of data set addresses stored in the
5 disk file are ordered;

6 (b) downloading at least one data set that includes an address of a referred data
7 set;

8 (c) identifying the address of the referred data set;

9 (d) generating a representation of the identified address;

10 (e) determining whether the representation is stored in the buffer, and whether the
11 disk file is empty;

12 (f) when the representation is not stored in the buffer and the disk file is empty,
13 scheduling the corresponding data set for downloading; and
14 (g) when the representation is not stored in the buffer and the disk file is not
15 empty, storing the representation in the buffer and delaying scheduling of the corresponding
16 data set for downloading until it is determined that the representation has not been previously
17 stored in the disk file.

1 23. A computer program product for use in conjunction with a computer system, the
2 computer program product comprising a computer readable storage medium and a computer
3 program mechanism embedded therein, the computer program mechanism comprising:
4 a first disk file and a buffer, for storing representations of data set addresses;
5 a main web crawler module for downloading and processing data sets stored on a
6 plurality of host computers, the main web crawler module identifying addresses of the one or
7 more referred data sets in the downloaded data sets; and
8 an address filtering module for processing a specified one of the identified addresses;
9 the address filtering module including instructions for:
10 generating a representation of the identified address;
11 determining whether the representation is stored in the buffer, and when this
12 determination is negative storing the representation in the buffer; and
13 determining whether the buffer has reached a predefined full condition, and
14 when this determination is positive, ordering the contents of the buffer and then performing
15 an ordered merge of contents of the buffer into the contents of the first disk file.

1 24. The computer program product of claim 23, wherein the address filtering module
2 further includes instructions for storing the identified address in the buffer after determining
3 that the representation is not stored in the buffer.

1 25. The computer program product of claim 23, wherein the address filtering module
2 further includes instructions for:
3 storing the identified address in a second disk file after determining that the
4 representation is not stored in the buffer; and

5 storing with each representation in the buffer a pointer to the corresponding address
6 stored in the second disk file; and
7 during the ordering of the contents of the buffer, keeping with each representation in
8 the buffer its pointer to the corresponding address in the second disk file.

1 26. The computer program product of claim 23, wherein
2 the first disk file is a sparse disk file divided into portions, each portion having a
3 starting address and contents comprising an ordered list of representations of data addresses;
4 and
5 the address filtering module includes instructions for performing the ordered merge of
6 the ordered contents of the buffer with the contents of the sparse disk file by obtaining a
7 starting address for a sub-file of the sparse disk file, the portion corresponding to one of the
8 representations in the buffer, and performing an ordered merge of a subset of the
9 representations in the buffer, starting at the one representation, into the contents of the
10 portion.

1 27. The computer program product of claim 23, wherein
2 the first disk file is a sparse disk file having empty entries interspersed among entries
3 storing said representations of data addresses; and
4 the address filtering module includes instructions for performing the ordered merge of
5 the ordered contents of the buffer with the contents of the sparse disk file by obtaining a
6 starting address corresponding to each respective representations in the buffer, and
7 sequentially scanning the first disk file, starting at the starting address, until the first of (A) a
8 representation matching the respective representation is found and (B) one of the empty
9 entries is found, and when an empty entry is found storing the respective representation in the
10 empty entry.

1 28. The computer program product of claim 23 wherein the representation of the
2 identified address comprises a checksum of at least a portion of the identified address.

1 29. The computer program product of claim 23, wherein the address filtering module
2 further includes instructions for first determining whether the representation is stored in a
3 cache, and when the first determination is positive, skipping the determination of whether the
4 representation is stored in the buffer.

1 30. The computer program product of claim 23, wherein the address filtering module
2 further includes instructions for:

3 determining whether the first disk file is empty and whether the representation is
4 stored in the buffer; and

5 if the first disk file is empty and the representation is not stored in the buffer, storing
6 the representation in the buffer and scheduling the corresponding data set for downloading.

1 31. A computer program product for use in conjunction with a computer system, the
2 computer program product comprising a computer readable storage medium and a computer
3 program mechanism embedded therein, the computer program mechanism comprising:

4 a first disk file, a first buffer, and a second buffer, for storing representations of data
5 set addresses;

6 a main web crawler module for downloading and processing data sets stored on a
7 plurality of host computers, the main web crawler module identifying addresses of the one or
8 more referred data sets in the downloaded data sets; and

9 an address filtering module for processing a specified one of the identified addresses;
10 the address filtering module including instructions for:

11 identifying one of the first and second buffers as a current buffer;
12 generating a representation of the identified address;
13 determining whether the representation is stored in the current buffer, and
14 when this determination is negative, storing the representation in the current buffer; and
15 determining whether the current buffer has reached a predefined full condition,
16 and when this determination is positive, selecting the other buffer as the current buffer,
17 wherein the previously current buffer is identified as a non-current buffer, ordering the
18 contents of the non-current buffer and then performing an ordered merge of the contents of
19 the non-current buffer into the contents of the first disk file.

1 32. The computer program product of claim 31, wherein the address filtering module
2 further includes instructions for storing the identified address in the current buffer after
3 determining that the representation is not stored in the current buffer.

1 33. The computer program product of claim 31, wherein the address filtering module
2 further includes instructions for:
3 storing the identified address in a second disk file after determining that the
4 representation is not stored in the current buffer;
5 storing with each representation in the current buffer a pointer to the corresponding
6 address stored in the second disk file; and
7 during the ordering of the contents of the non-current buffer, keeping with each
8 representation in the non-current buffer its pointer to the corresponding address in the second
9 disk file.

1 34. The computer program product of claim 31, wherein
2 the first disk file is a sparse disk file divided into sub-files, each sub-file having a
3 starting address and contents comprising an ordered list of representations of data addresses;
4 and
5 the instructions for performing the ordered merge including instructions for obtaining
6 a starting address for a sub-file of the first disk file, the sub-file corresponding to one of the
7 representations in the buffer, and performing an ordered merge of a subset of the
8 representations in the non-current buffer, starting at the one representation, into the contents
9 of the sub-file.

10 35. The computer program product of claim 31, wherein
11 the first disk file is a sparse disk file having empty entries interspersed among entries
12 storing said representations of data addresses; and
13 the address filtering module includes instructions for performing the ordered merge of
14 the ordered contents of the buffer with the contents of the sparse disk file by obtaining a
15 starting address corresponding to each respective representations in the buffer, and

16 sequentially scanning the first disk file, starting at the starting address, until the first of (A) a
17 representation matching the respective representation is found and (B) one of the empty
18 entries is found, and when an empty entry is found storing the respective representation in the
19 empty entry.

1 36. The computer program product of claim 31 wherein the representation of the
2 identified address comprises a checksum of at least a portion of the identified address.

1 37. The computer program product of claim 31, wherein the address filtering module
2 further includes instructions for:

3 determining whether the first disk file is empty and whether the representation is
4 stored in the current buffer; and

5 if the first disk file is empty and the representation is not stored in the current buffer,
6 storing the representation in the current buffer and scheduling the corresponding data set for
7 downloading.

1 38. A web crawler for downloading data set addresses from among a plurality of host
2 computers, comprising:

3 a first disk file and a buffer, for storing representations of data set addresses;

4 a main web crawler module for downloading and processing data sets stored on a
5 plurality of host computers, the main web crawler module identifying addresses of the one or
6 more referred data sets in the downloaded data sets; and

7 an address filtering module for processing a specified one of the identified addresses;
8 the address filtering module including instructions for:

9 generating a representation of the identified address;

10 determining whether the representation is stored in the buffer, and when this
11 determination is negative storing the representation in the buffer; and

12 determining whether the buffer has reached a predefined full condition, and
13 when this determination is positive, ordering the contents of the buffer and then performing
14 an ordered merge of the contents of the buffer into the contents of the first disk file.

1 39. The web crawler of claim 38, wherein the address filtering module further includes
2 instructions for storing the identified address in the buffer following a determination that the
3 representation is not stored in the buffer.

1 40. The web crawler of claim 38, wherein the address filtering module further includes
2 instructions for:

3 storing the identified address in a second disk file after determining that the
4 representation is not stored in the buffer; and

5 storing with each representation in the buffer a pointer to the corresponding address
6 stored in the second disk file; and

7 during the ordering of the contents of the buffer, keeping with each representation in
8 the buffer its pointer to the corresponding address in the second disk file.

1 41. The web crawler of claim 38 wherein

2 the first disk file is a sparse disk file divided into portions, each portion having a
3 starting address and contents comprising an ordered list of representations of data addresses;
4 and

5 the address filtering module further includes instructions for:

6 obtaining, from an index, a starting address for a portion in the sparse disk file
7 corresponding to one of the representations stored in the buffer; and

8 performing an ordered merge of a subset of the representations stored in the
9 buffer, starting at the representation for which the starting address was obtained, into the
10 contents of the corresponding portion.

1 42. The web crawler of claim 38 wherein

2 the first disk file is a sparse disk file having empty entries interspersed among entries
3 storing said representations of data addresses; and

4 the address filtering module includes instructions for performing the ordered merge of
5 the ordered contents of the buffer with the contents of the sparse disk file by obtaining a
6 starting address corresponding to each respective representations in the buffer, and
7 sequentially scanning the first disk file, starting at the starting address, until the first of (A) a

8 representation matching the respective representation is found and (B) one of the empty
9 entries is found, and when an empty entry is found storing the respective representation in the
10 empty entry.

1 43. The web crawler of claim 38 wherein the representation of the identified address
2 comprises a checksum of at least a portion of the identified address.

1 44. The web crawler of claim 38 wherein the address filtering module further includes
2 instructions for:

3 determining whether the representation is stored in a cache before determining
4 whether the representation is stored in the buffer, and when this determination is negative,
5 determining whether the representation is stored in the buffer;

6 when the second determination is negative, storing the representation in the buffer;

7 when the first determination is negative, and predefined other criteria are met, storing
8 the representation in the cache; and

9 when the cache has reached a predefined full condition, evicting a stored
10 representation from the cache in accordance with an eviction policy.

1 45. The web crawler of claim 38 wherein the address filtering module further includes
2 instructions for determining whether the first disk file is empty and whether the
3 representation is stored in the buffer, and if the first disk file is empty and the representation
4 is not stored in the buffer, storing the representation in the buffer and scheduling the
5 corresponding data set for downloading.

1 46. A web crawler for downloading data set addresses from among a plurality of host
2 computers, comprising:

3 a first disk file, a first buffer and a second buffer, for storing representations of data
4 set addresses;

5 a main web crawler module for downloading and processing data sets stored on a
6 plurality of host computers, the main web crawler module identifying addresses of the one or
7 more referred data sets in the downloaded data sets; and

an address filtering module for processing a specified one of the identified addresses;
the address filtering module including instructions for:
identifying one of the first and second buffers as a current buffer;
generating a representation of the identified address;
determining whether the representation is stored in the current buffer, and
when this determination is negative, storing the representation in the current buffer; and
determining whether the current buffer has reached a predefined full condition,
and when this determination is positive, selecting the other buffer as the current buffer,
wherein the previously current buffer is identified as a non-current buffer, ordering the
contents of the non-current buffer and then performing an ordered merge of the contents of
the non-current buffer into the contents of the first disk file.

47. The web crawler of claim 46, wherein the address filtering module further includes
instructions for storing the identified address in the current buffer after determining that the
representation is not stored in the current buffer.

48. The web crawler of claim 46, wherein the address filtering module further includes
instructions for:
storing the identified address in a second disk file after determining that the
representation is not stored in the current buffer;
storing with each representation in the current buffer a pointer to the corresponding
address stored in the second disk file; and
during the ordering of the contents of the non-current buffer, keeping with each
representation in the non-current buffer its pointer to the corresponding address in the second
disk file.

49. The web crawler of claim 46, wherein
the first disk file is a sparse disk file divided into sub-files, each sub-file having a
starting address and contents comprising an ordered list of representations of data addresses;
and

5 the instructions for performing the ordered merge including instructions for obtaining
6 a starting address for a sub-file of the first disk file, the sub-file corresponding to one of the
7 representations in the buffer, and performing an ordered merge of a subset of the
8 representations in the non-current buffer, starting at the one representation, into the contents
9 of the sub-file.

1 50. The web crawler of claim 46 wherein

2 the first disk file is a sparse disk file having empty entries interspersed among entries
3 storing said representations of data addresses; and

4 the address filtering module includes instructions for performing the ordered merge of
5 the ordered contents of the buffer with the contents of the sparse disk file by obtaining a
6 starting address corresponding to each respective representations in the buffer, and
7 sequentially scanning the first disk file, starting at the starting address, until the first of (A) a
8 representation matching the respective representation is found and (B) one of the empty
9 entries is found, and when an empty entry is found storing the respective representation in the
10 empty entry.

1 51. The web crawler of claim 46 wherein the representation of the identified address
2 comprises a checksum of at least a portion of the identified address.

1 52. The web crawler of claim 46, wherein the address filtering module further includes
2 instructions for:

3 determining whether the first disk file is empty and whether the representation is
4 stored in the current buffer; and

5 when the first disk file is empty and the representation is not stored in the current
6 buffer, storing the representation in the current buffer and scheduling the corresponding data
7 set for downloading.